

# Scalable Factorized Hierarchical Variational Autoencoder Training

Wei-Ning Hsu, James Glass

Computer Science and Artificial Intelligence Laboratory  
Massachusetts Institute of Technology  
Cambridge, MA 02139, USA  
{wnhsu, glass}@mit.edu

## Abstract

Deep generative models have achieved great success in unsupervised learning with the ability to capture complex nonlinear relationships between latent generating factors and observations. Among them, a factorized hierarchical variational autoencoder (FHVAE) is a variational inference-based model that formulates a hierarchical generative process for sequential data. Specifically, an FHVAE model can learn disentangled and interpretable representations, which have been proven useful for numerous speech applications, such as speaker verification, robust speech recognition, and voice conversion. However, as we will elaborate in this paper, the training algorithm proposed in the original paper is not scalable to datasets of thousands of hours, which makes this model less applicable on a larger scale. After identifying limitations in terms of runtime, memory, and hyperparameter optimization, we propose a hierarchical sampling training algorithm to address all three issues. Our proposed method is evaluated comprehensively on a wide variety of datasets, ranging from 3 to 1,000 hours and involving different types of generating factors, such as recording conditions and noise types. In addition, we also present a new visualization method for qualitatively evaluating the performance with respect to interpretability and disentanglement. Models trained with our proposed algorithm demonstrate the desired characteristics on all the datasets.

**Index Terms:** unsupervised learning, speech representation learning, factorized hierarchical variational autoencoder

## 1. Introduction

Unsupervised learning can leverage large amounts of unlabeled data to discover latent generating factors that often lie on a lower dimensional manifold compared to the raw data. A learned latent representation from speech can be useful for many downstream applications, such as speaker verification [1], automatic speech recognition [2], and linguistic unit discovery [3, 4]. A factorized hierarchical variational autoencoder (FHVAE) [5] is a variational inference-based deep generative model that learns interpretable and disentangled latent representation from sequential data without supervision by modeling a hierarchical generative process. In particular, it has been demonstrated that an FHVAE trained on speech data learns to encode sequence-level generating factors, such as speaker and channel condition, into one set of latent variables, while encoding segment-level generating factors, such as phonetic content, into another set of latent variables. The ability to disentangle latent factors has been beneficial to a wide range of tasks, including domain adaptation [6], conditional data augmentation [7], and voice conversion [5].

However, the original FHVAE training algorithm proposed in [5] does not easily scale to datasets of over hundreds of thousands of utterances, making it less applicable to real world settings, where an unlabeled dataset of such size is common. This limitation is mainly due to the following issues: (1) the inference model of the sequence-level latent variable, and (2) the design of the discriminative objective. To be more specific, the original training algorithm reduces the complexity of inferring sequence-level latent variables by maintaining a lookup table, whose number of entries equals the number of training sequences. In addition, the discriminative object, which encourages disentanglement, requires computing a partition function that sums over a function of each entry in that lookup table. The two facts combined lead to significant scalability issues.

In this paper, we propose a hierarchical sampling algorithm to address these issues. In addition, a new method for qualitatively evaluating disentanglement performance based on a t-Distribution Stochastic

Neighbor Embedding [8] is also presented. The proposed training algorithm is evaluated on a wide variety of datasets, ranging from 3 to 1,000 hours and involving many different types of generating factors, such as recording conditions and noise types. Experimental results verify that the proposed algorithm is effective on all sizes of datasets and universally achieves desirable disentanglement performance. All codes are available on-line.<sup>1</sup>

## 2. Limitations of Original FHVAE Training

In this section, we briefly review the factorized hierarchical variational autoencoder (FHVAE) model and discuss the scalability issues of the original training objective.

### 2.1. Factorized Hierarchical Variational Autoencoders

An FHVAE [5] is a variant of VAEs that models a generative process of sequential data with a hierarchical graphical model. Specifically, generation of a segment  $\mathbf{x}$ , which is a sub-sequence of a sequence  $\mathbf{X} = \{\mathbf{x}^{(n)}\}_{n=1}^N$ , is conditioned on a pair of segment-level latent variables,  $\mathbf{z}_1$  and  $\mathbf{z}_2$ , referred to as the *latent segment variable* and the *latent sequence variable* respectively. While  $\mathbf{z}_1$  is generated from a global prior, similar to those latent variables in a vanilla VAE,  $\mathbf{z}_2$  is generated from a sequence-dependent prior that is conditioned on a sequence-level latent variable,  $\mu_2$ , named the *s-vector*. The joint probability of a sequence is formulated as:  $p(\mu_2) \prod_{n=1}^N p(\mathbf{x}^{(n)} | \mathbf{z}_1^{(n)}, \mathbf{z}_2^{(n)}) p(\mathbf{z}_1^{(n)}) p(\mathbf{z}_2^{(n)} | \mu_2)$ . With this formulation, an FHVAE learns to encode generating factors that are consistent within segments drawn from the same sequence into  $\mathbf{z}_2$ . In contrast,  $\mathbf{z}_1$  captures the residual generating factors that changes between segments.

Since computing the true posteriors of  $\mathbf{Z}_1 = \{\mathbf{z}_1^{(n)}\}_{n=1}^N$ ,  $\mathbf{Z}_2 = \{\mathbf{z}_2^{(n)}\}_{n=1}^N$ , and  $\mu_2$  are intractable, an FHVAE introduces an inference model,  $q(\mathbf{Z}_1, \mathbf{Z}_2, \mu_2 | \mathbf{X})$ , and factorizes it as:  $q(\mu_2 | \mathbf{X}) \prod_{n=1}^N q(\mathbf{z}_1^{(n)} | \mathbf{x}^{(n)}) q(\mathbf{z}_2^{(n)} | \mathbf{x}^{(n)})$ . We summarize in Table 1 the family of distributions an FHVAE adopts for the generative model and the inference model. All the functions,  $f_{\mu_{\mathbf{x}}}(\cdot, \cdot)$ ,  $f_{\sigma_{\mathbf{x}}}(\cdot, \cdot)$ ,  $g_{\mu_{\mathbf{z}_1}}(\cdot, \cdot)$ ,  $g_{\sigma_{\mathbf{z}_1}}(\cdot, \cdot)$ , and  $g_{\mu_{\mathbf{z}_2}}(\cdot, \cdot)$ ,  $g_{\sigma_{\mathbf{z}_2}}(\cdot, \cdot)$ , are neural networks that parameterize mean and variance of Gaussian distributions.

Table 1: Family of distributions adopted for FHVAE generative and inference models.

generative model	
$p(\mu_2)$	$\mathcal{N}(\mathbf{0}, \mathbf{I})$
$p(\mathbf{z}_1)$	$\mathcal{N}(\mathbf{0}, \mathbf{I})$
$p(\mathbf{z}_2   \mu_2)$	$\mathcal{N}(\mu_2, \sigma_{\mathbf{z}_2}^2 \mathbf{I})$
$p(\mathbf{x}   \mathbf{z}_1, \mathbf{z}_2)$	$\mathcal{N}(f_{\mu_{\mathbf{x}}}(\mathbf{z}_1, \mathbf{z}_2), \text{diag}(f_{\sigma_{\mathbf{x}}}^2(\mathbf{z}_1, \mathbf{z}_2)))$
inference model	
$q(\mu_2   \mathbf{X})$	$\mathcal{N}(\sum_{n=1}^N g_{\mu_{\mathbf{z}_2}}(\mathbf{x}^{(n)}) / (N + \sigma_{\mathbf{z}_2}^2), \mathbf{I})$
$q(\mathbf{z}_1   \mathbf{x}, \mathbf{z}_2)$	$\mathcal{N}(g_{\mu_{\mathbf{z}_1}}(\mathbf{x}, \mathbf{z}_2), \text{diag}(g_{\sigma_{\mathbf{z}_1}}^2(\mathbf{x}, \mathbf{z}_2)))$
$q(\mathbf{z}_2   \mathbf{x})$	$\mathcal{N}(g_{\mu_{\mathbf{z}_2}}(\mathbf{x}), \text{diag}(g_{\sigma_{\mathbf{z}_2}}^2(\mathbf{x})))$

<sup>1</sup><https://github.com/wnhsu/ScalableFHVAE>

## 2.2. Original FHVAE Training

In the variational inference framework, since the marginal likelihood of observed data is intractable, we optimize a lower bound of it instead, named variational lower bound,  $\mathcal{L}(p, q; \mathbf{X})$ . We can derive a sequence variational lower bound of  $\mathbf{X}$  based on Table 1. However, this lower bound can only be optimized at the sequence level, because inferring of  $\mu_2$  depends on an entire sequence, and would become infeasible if  $\mathbf{X}$  is extremely long.

In [5] the authors proposed replacing the maximum a posterior (MAP) estimation of  $\mu_2$ 's posterior mean for training sequences with a lookup table  $h_{\mu_2}(i)$ , where  $i$  indexes training sequences. In other words, the inference model for  $\mu_2$  becomes  $q(\mu_2 | \mathbf{X}^{(i)}) = \mathcal{N}(h_{\mu_2}(i), \mathbf{I})$ . Therefore, the lower bound can be re-written as:

$$\mathcal{L}(p, q; \mathbf{X}^{(i)}) = \sum_{n=1}^{N(i)} \mathcal{L}(p, q; \mathbf{x}^{(i,n)} | h_{\mu_2}(i)) + \log p(h_{\mu_2}(i)) \quad (1)$$

$$\begin{aligned} \mathcal{L}(p, q; \mathbf{x}^{(i,n)} | h_{\mu_2}(i)) &= \mathbb{E}_{q(\mathbf{z}_1, \mathbf{z}_2 | \mathbf{x}^{(i,n)})} [\log p(\mathbf{x}^{(i,n)} | \mathbf{z}_1, \mathbf{z}_2)] \\ &\quad - \mathbb{E}_{q(\mathbf{z}_2 | \mathbf{x}^{(i,n)})} [D_{KL}(q(\mathbf{z}_1 | \mathbf{x}^{(i,n)}, \mathbf{z}_2) || p(\mathbf{z}_1))] \\ &\quad - D_{KL}(q(\mathbf{z}_2 | \mathbf{x}^{(i,n)}) || p(\mathbf{z}_2 | h_{\mu_2}(i))). \end{aligned} \quad (2)$$

We can now sample a batch at the segment level to optimize the following *segment variational lower bound*:

$$\mathcal{L}(p, q; \mathbf{x}^{(i,n)}) = \mathcal{L}(p, q; \mathbf{x}^{(i,n)} | h_{\mu_2}(i)) + \frac{1}{N(i)} \log p(h_{\mu_2}(i)). \quad (3)$$

Furthermore, to obtain meaningful disentanglement between  $\mathbf{z}_1$  and  $\mathbf{z}_2$ , it is not desirable to have constant  $\mu_2$  for all sequences, which would result in interchangeable  $\mathbf{z}_1$  and  $\mathbf{z}_2$ . To avoid such condition, the following objective is added to encourage  $\mu_2$  to be discriminative between sequences:

$$\log p(i | \bar{\mathbf{z}}_2^{(i,n)}) := \log \frac{p(\bar{\mathbf{z}}_2^{(i,n)} | \bar{\mu}_2^{(i)})}{\sum_{j=1}^M p(\bar{\mathbf{z}}_2^{(i,n)} | \bar{\mu}_2^{(j)})}, \quad (4)$$

where  $M$  is the total number of training sequences,  $\bar{\mathbf{z}}_2^{(i,n)}$  denotes the posterior mean of  $\mathbf{z}_2$ ,  $g_{\mu_2}(\mathbf{x}^{(i,n)})$ , and  $\bar{\mu}_2^{(i)}$  denotes the posterior mean of  $\mu_2$ ,  $h_{\mu_2}(i)$ . This additional discriminative objective encourages  $\mathbf{z}_2$  from the  $i$ -th sequence to be not only close to  $\mu_2$  of the  $i$ -th sequence, but also far from  $\mu_2$  of other sequences. Combining this discriminative objective and the segment variational lower bound with a weighting parameter,  $\alpha$ , the objective function that an FHVAE maximizes then becomes:

$$\mathcal{L}^{dis}(p, q; \mathbf{x}^{(i,n)}) = \mathcal{L}(p, q; \mathbf{x}^{(i,n)}) + \alpha \log p(i | \bar{\mathbf{z}}_2^{(i,n)}), \quad (5)$$

referred to as the *discriminative segment variational lower bound*.

## 2.3. Scalability Issues

The original FHVAE training addressed the scalability issue with respect to sequence length by decomposing a sequence variational lower bound into sum of segmental variational lower bounds over segments. However, here we will show that this training objective is not scalable with respect to the number of training sequences.

First of all, the original FHVAE training maintains a lookup table,  $h_{\mu_2}(\cdot)$ , that stores the posterior mean of  $\mu_2$  for each training sequence. The size of this table grows proportionally to the number of training sequences. However, this is less problematic, unless the number of our training sequences is on the order of  $10^8$ , which can make the lookup table size grow to 10 GB (suppose  $\mathbf{z}_2$  are 32-dimensional 32-bit floating point vectors). The main limitation here is the computation of the discriminative objective. The denominator,  $\sum_{j=1}^M p(\bar{\mathbf{z}}_2^{(i,n)} | \bar{\mu}_2^{(j)})$ , marginalizes over a function of posterior mean of  $\mu_2$  of all training sequences, which increases the computation time proportionally to the number of sequences. Furthermore, when computing the gradient given a batch of training segments, we need to maintain a tensor of size  $(bs, |\theta|)$ , where  $bs$  is the batch size, and  $|\theta|$  is the total number of trainable parameters involved in the computation of the objective function.

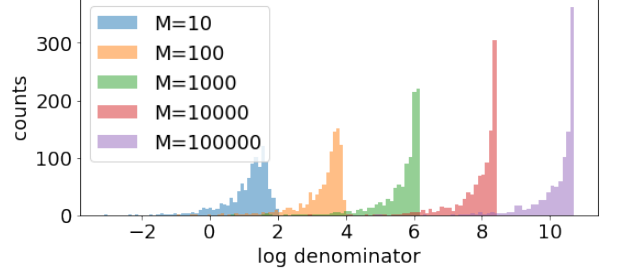


Figure 1: Histogram of  $\log \sum_{i=1}^M p(\mathbf{z}_2^{(i,n)} | \bar{\mu}_2^{(j)})$  with respect to different  $M \in \{10^1, 10^2, 10^3, 10^4, 10^5\}$ . Distributions shift by roughly a constant when  $M$  increases by 10 times, implying the denominator scales proportionally to  $M$ .

Since the computation of the discriminative objective involves the entire lookup table,  $h_{\mu_2}(\cdot)$ , the gradient tensor is of size at least  $bs$  times larger than the lookup table. With a batch size of 256, a dataset with  $10^5$  sequences can exhaust the memory of a single GPU during training.

In addition, from the hyperparameter optimization point of view, it can be observed that the distribution of the denominator,  $\sum_{i=1}^M p(\mathbf{z}_2^{(i,n)} | \bar{\mu}_2^{(j)})$ , also changes with respect to the number of training sequences. Specifically, the expected value of this terms scales proportionally to  $M$  asymptotically, as shown in Figure 1. Such behavior is not desirable, because the  $\alpha$  parameter that balances the variational lower bound and the discriminative objective would need to be adjusted according to  $M$ .

## 3. Training with Hierarchical Sampling

In order to utilize the discriminative objective, while eliminating the memory, computation, and optimization issue induced by a large training set, we need to control the size of the lookup table as well as the denominator summation in the discriminative objective. Both of these can be achieved jointly with a hierarchical sampling algorithm.

Given a dataset of  $M$  training sequences, we maintain a lookup table of only  $K$  entries, where  $K$  is a dataset independent hyperparameter, and optimize an FHVAE model with the following procedure: (1) A batch of  $K$  sequences,  $\{\tilde{\mathbf{X}}^{(k)}\}_{k=1}^K$ , is sampled from the entire training set, where  $k \in [1, K]$  indexes the sampled training sequences. (2) The  $k$ -th entry of the lookup table,  $h_{\mu_2}(k)$ , is updated with the MAP estimation of the  $k$ -th sequence in the batch,  $\sum_{n=1}^{N(k)} g_{\mu_2}(\tilde{\mathbf{x}}^{(k,n)}) / (N + \sigma_{\mathbf{z}_2}^2)$ . (3)  $N_{Bseg}$  batches of segments are drawn from the  $K$  sequences. (4) Each segment batch is used iteratively to estimate the discriminative segmental variational lower bound for optimizing the parameters of  $f_*$ ,  $g_*$ , and  $h_{\mu_2}$  as before. The only difference is that the denominator of the discriminative object now sums over the  $K$  sampled training sequences, instead of the entire set of  $M$  training sequences. (5) repeat the above steps until convergence. We list the pseudo code in Algorithm 1.

We refer to the proposed algorithm as a hierarchical sampling algorithm, because we first sample at the sequence level, and then at the segment level, in order to reduce the effective number of training sequences from  $M$  to  $K$  for each segment batch. Compared with the proposed algorithm, the original training algorithm can be regarded as using a “flat” sampling algorithm, where we sample segments from the entire pool, so it is therefore necessary to maintain a lookup table of  $M$  entries. The proposed algorithm introduces overhead associated with updating the lookup table whenever a new batch of sequences is sampled. However, this cost can be amortized by increasing the number of segment batches,  $N_{Bseg}$ , for each batch of sequences.

## 4. Experimental Setup

We evaluate our training algorithm on a wide variety of datasets, ranging from 3 to 1,000 hours, including both clean and noisy, close-talking and distant speech. In this section, we describe the datasets, and introduce FHVAE models and their training configurations.

---

**Algorithm 1** Training with Hierarchical Sampling
 

---

**Input:**  $\{\mathbf{X}^{(i)}\}_{i=1}^M$ : training set;  $K$ : sequence batch size;  $bs$ : segment batch size;  $NB_{seg}$ : number of segment batches;  $f_*/g_*$ : decoders/encoders;  $h_{\mu_{\mu_2}}$ : lookup table of  $K$  entries;  $Opt$ : gradient descent-based optimizer

```

1: while not converged do
2:   sample a batch of  $K$  training sequences,  $\{\tilde{\mathbf{X}}^{(k)}\}_{k=1}^K$ 
3:   for  $k = 1 \dots K$  do
4:      $h_{\mu_{\mu_2}}(k) \leftarrow \sum_{n=1}^{N^{(k)}} g_{\mu_{\mu_2}}(\tilde{\mathbf{x}}^{(k,n)}) / (N + \sigma_{\mathbf{z}_2}^2)$ 
5:   end for
6:   for  $1 \dots NB_{seg}$  do
7:     sample segments  $\{\tilde{\mathbf{x}}^{(k_b, n_b)}\}_{b=1}^{bs}$  from  $\{\tilde{\mathbf{X}}^{(k)}\}_{k=1}^K$ 
8:      $loss_{dis}(b) \leftarrow \log \frac{p(g_{\mu_{\mu_2}}(\tilde{\mathbf{x}}^{(k_b, n_b)}) | h_{\mu_{\mu_2}}(k_b))}{\sum_{k=1}^K p(g_{\mu_{\mu_2}}(\tilde{\mathbf{x}}^{(k_b, n_b)}) | h_{\mu_{\mu_2}}(k))}$ 
9:      $loss_{gen}(b) \leftarrow \mathcal{L}(p, q; \tilde{\mathbf{x}}^{(k_b, n_b)})$ 
10:     $loss \leftarrow -\sum_{b=1}^{bs} (loss_{gen}(b) + \alpha * loss_{dis}(b)) / bs$ 
11:     $f_*, g_*, h_{\mu_{\mu_2}} \leftarrow Opt(loss, \{f_*, g_*, h_{\mu_{\mu_2}}\})$ 
12:   end for
13: end while
14: return  $f_*, g_*$ 

```

---

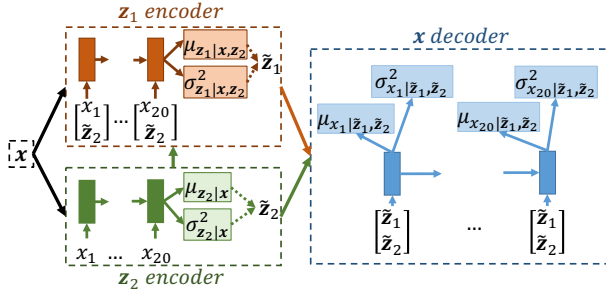


Figure 2: The proposed FHVAE architecture consists of two encoders (orange and green) and one decoder (blue).  $\mathbf{x} = [x_1, \dots, x_{20}]$  is a segment of 20 frames. Dotted lines in the encoders denote sampling from parametric distributions.

#### 4.1. Datasets

Four different corpora are used for our experiments, which are TIMIT [9], Aurora-4 [10], AMI [11], and LibriSpeech [12]. TIMIT contains broadband 16kHz recordings of phonetically-balanced read speech. A total of 3,696 utterances (3 hours) are presented in the training partition based on the Kaldi [13] recipe, where *sa* utterances are excluded. In addition, manually labeled time-aligned phonetic transcripts are available, which we use to study the disentanglement performance between phonetic and speaker information achieved by FHVAE models.

Aurora-4 is another broadband corpus designed for noisy speech recognition tasks, based on the Wall Street Journal corpus (WSJ0) [14]. Six different types of noise are artificially mixed with clean read speech of two different microphone types, amounting to a total of 4,620 utterances (10 hours) for the development set. Following the training setup in [5], we train our FHVAE models on the development set, because the noise and channel information on the training set is not available.

The AMI corpus consists of 100 hours of meeting recordings, recorded in three different meeting rooms with different acoustic properties, and with multiple attendants. Multiple microphones are used for each session, including individual headset microphones (IHM), and far-field microphone arrays. IHM and single distant microphone (SDM) recordings from the training set are mixed to form a training set for the FHVAE models, including over 200,000 utterances according to the segmentation provided in the corpus.

The largest corpus we evaluate on is the LibriSpeech corpus, which contains 1,000 hours of read speech sampled at 16kHz. This corpus is based on the LibriVox’s project, where world-wide volunteers record public domain texts to create free public domain audio books.

#### 4.2. Training and Model Configurations

Speech segments of 20 frames, represented with 80-dimensional log Mel-scale filter bank coefficients (FBank), are used as inputs to FHVAE models. We denote each segment with  $\mathbf{x} = [x_1, \dots, x_{20}]$ . The variance of  $\mathbf{z}_2$ ’s prior is set to  $\sigma_{\mathbf{z}_2}^2 = 0.25$ , and the dimension of  $\mathbf{z}_1$  and  $\mathbf{z}_2$  are both 32. Figure 2 illustrates the detailed encoder/decoder architectures of the proposed FHVAE model. The conditional mean and variance predictor for each variable (i.e.,  $\mathbf{z}_1$ ,  $\mathbf{z}_2$ , and  $\mathbf{x}$ ) shares a common stacked LSTM pre-network, followed by two different single-layer affine transform networks,  $\mu_*$  and  $\sigma_*^2$ , predicting the conditional mean and variance respectively. Specifically, a stacked LSTM with 2 layers and 256 memory cells are used for all three pre-networks, illustrated in Figure 2 with blocks filled with dark colors. Affine transform networks of  $\mathbf{z}_1$  and  $\mathbf{z}_2$  encoders take as input the output from the last time step of both layers, which sums to 512 dimension. As for the  $\mathbf{x}$  decoder, the affine transform network takes as input the LSTM output of the last layer from each time step  $t$ , and predicts the probability distribution of the corresponding frame  $p(x_t | \mathbf{z}_1, \mathbf{z}_2)$ . The same sampled  $\tilde{\mathbf{z}}_1$  and  $\tilde{\mathbf{z}}_2$  from the posterior distributions are concatenated and used as input for the LSTM decoder at each step. Sampling is done by introducing auxiliary input variables for the reparameterization trick [15], in order to keep the entire network differentiable with respect to the objective.

FHVAE models are trained to optimize the discriminative segment variational lower bound with  $\alpha = 10$ . We set sequence batch size  $K = 2000$  for TIMIT and Aurora-4, and  $K = 5000$  for the others. Adam [16] with  $\beta_1 = 0.95$  and  $\beta_2 = 0.999$  is used to optimize all models. Tensorflow [17] is used for implementation. Training is done for 500,000 steps, terminating early if the segmental variational lower bound on a held-out validation set is not improved for 50,000 steps.

## 5. Results and Discussions

#### 5.1. Time and Memory Complexity

One feature of our proposed training algorithm is to control memory complexity. We found that a training set with over 100,000 sequences would exhaust a single 8GB GPU memory when using the original training algorithm. Hence, it was not feasible for the AMI and the LibriSpeech corpus, while the proposed algorithm does not suffer from the same problem. Another feature of hierarchical sampling is to control the time complexity of computing the discriminative loss. To study how sequence batch size affects the optimization step (line 8 in Algorithm 1), we evaluate the processing time of that step by varying  $K$  from 20 to 20,000 and show the results in Table 2.

We can observe that when  $K \leq 2000$ , the time complexity of computing the discriminative loss is fractional compared to computing the variational lower bound. However, when  $K > 2000$ , the increased computation time grows proportional to the sequence batch size, so that computation of the discriminative loss starts to dominate the time complexity. In practice, given a new encoder/decoder architecture, we can investigate the computation overhead resulting from the discriminative loss using such a method, and it is possible to determine some  $K$  that introduces negligible overhead for optimization.

Table 2: Processing time of the optimization step with different sequence batch size  $K$ .

$K$	10	100	1000	2000	5000	10000	20000
Time (ms)	84	84	86	87	103	147	230

#### 5.2. Evaluating Disentanglement Performance

To examine whether an FHVAE is successfully trained, we need to inspect its performance at disentangling sequence-level generating factors (e.g. speaker identity, noise condition, and channel condition) from segment-level generating factors (e.g. phonetic content) in the latent space. For quantitative evaluation, we reproduce the speaker verification experiments in [5]. The FHVAE model trained with hierarchical sampling achieves 1.64% equal error rate on TIMIT, matching the performance of the original training algorithm (1.34%). In the following sections, we proceed with two qualitative evaluation methods.



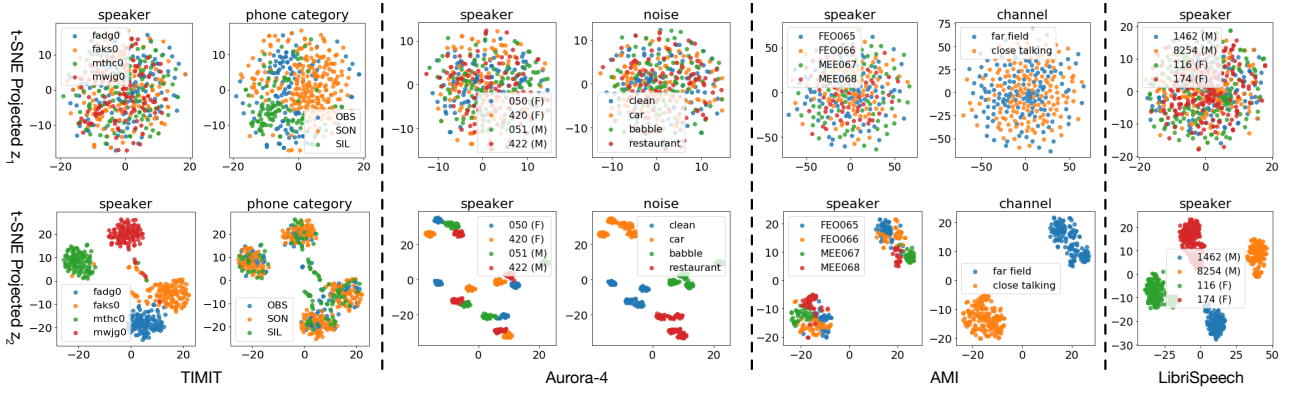


Figure 3: Scatter plots of t-SNE projected  $z_1$  and  $z_2$  with models trained on TIMIT/Aurora-4/AMI/LibriSpeech. Each point represents one segment. Different colors are used to code segments of different labels with respect to the generating factor shown at the title of each plot.

### 5.2.1. t-SNE Visualization of Latent Variables

We start with selecting a batch of labeled segments  $(x, y)$ , where  $y$  denotes the values of the associated generating factors, for example  $y = (\text{phone-id}, \text{speaker-id})$ . We then infer  $z_1$  and  $z_2$  of these segments, and project them separately to a two-dimensional space using t-Distributed Stochastic Neighbor Embedding (t-SNE) [8]. Each generating factor is used to color-code both projected  $z_1$  and  $z_2$ . Successful disentanglement would result in segments of the same sequence-level generating factors forming clusters in the projected  $z_2$  space but not in the projected  $z_1$  space, and vice versa.

For all four datasets, speaker label, a sequence-level generating factor, is available for each segment. Since time-aligned phonetic transcripts are available for TIMIT, it is also possible to derive phone labels, which is a segment-level generating factor. Following [18], we further reduce the 61 phonemes to three phonetic subsets: sonorant (SON), obstruent (OBS), and silence (SIL) for better color-coding. In addition, noise types can be obtained for Aurora-4, and microphone types can be obtained for AMI, which are both sequence-level generating factors.

Results of t-SNE projections for models trained on each dataset are shown in Figure 3, where each point represents one segment. It can be observed that in each of the projected  $z_2$  spaces, segments of the same sequence-level generating factors (speaker/noise/channel) always form clusters. When segments are generated conditioned on multiple sequence-level generating factors, as in Aurora-4 and AMI, the segments actually cluster hierarchically. In contrast, the distribution of projected  $z_1$ 's does not vary between different values of these generating factors, which implies that  $z_1$  does not contain much information about them. Opposite phenomenon can be observed from the phone category-coded plots, where segments belonging to the same phonetic subset cluster in the projected  $z_1$  space, but not in the projected  $z_2$  space. These results suggest that FHVAEs trained with hierarchical sampling can achieve desirable disentanglement for these conditions.

### 5.2.2. Reconstructing Re-combined Latent Variables

Given two segments,  $x^A$  and  $x^B$ , we sample a segment  $x^C \sim p(x|\tilde{z}_1^A, \tilde{z}_2^B)$ , where  $\tilde{z}_1^A$  is a sampled latent segment variable conditioned on  $x^A$ , and  $\tilde{z}_2^B$  is a sampled latent sequence variable conditioned on  $x^B$ . With a successfully trained FHVAE,  $x^C$  should exhibit the segment-level attributes of  $x^A$ , and the sequence-level attributes of  $x^B$ . Due to space limitations, we only show results of the model trained on the AMI corpus in Figure 4. Eight segments are sampled for  $x^A$ , as shown in the upper right corner of the figure. Among these segments, the four leftmost ones are close-talking while the rest are far-field, and the first, second, fifth, and sixth from left are female speakers while the rest are males. For  $x^B$ , three segments of different sequence-level generating factors are sampled, as shown on the left half of the figure. The segments used to infer  $\tilde{z}_2^B$  are highlighted in red boxes; we show the surrounding frames of those segments to better illustrate how sequence-level generating factors affect realization of observations.

Samples of  $x^C$  generated by re-combining latent variables are shown in the lower right corner of the figure. It can be clearly ob-

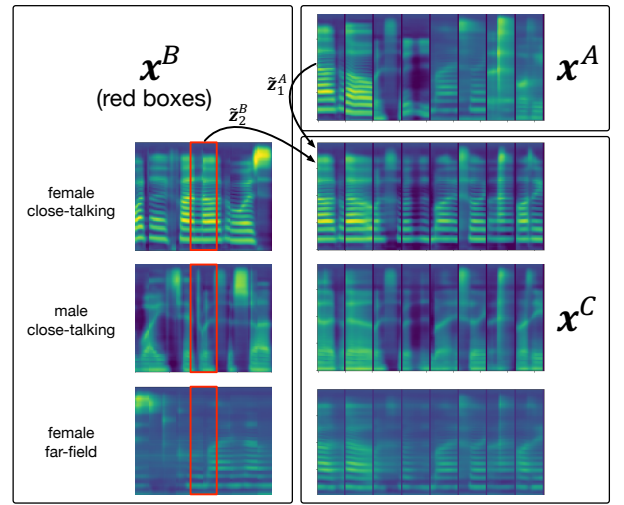


Figure 4: Results of decoding re-combined latent variables. A segment in the  $x^C$  block is generated conditioned on the latent segment variable of a segment in the block  $x^A$  of the same column, and conditioned on the latent sequence variable of a red-box highlighted segment in the block  $x^B$  of the same row.

served that  $x^C$  presents the same sequence-level generating factors as  $x^B$ ,<sup>2</sup> whose latent sequence variable  $x^C$  conditions on. Meanwhile, the phonetic content of  $x^C$  stays consistent with  $x^A$ ,<sup>3</sup> whose latent segment variable  $x^C$  conditions on. The clear differentiation of generating factors encoded in each sets of latent variables again corroborates the success of our proposed algorithm in training FHVAE models.

## 6. Conclusions and Future Work

In this paper, we discuss the scalability limitations of the original FHVAE training algorithm in terms of runtime, memory, and hyperparameter optimization, and propose a hierarchical sampling algorithm to address this problem. Comprehensive study on the memory and time complexity, as well as disentanglement performance verify the effectiveness of the proposed algorithm on all scales of datasets, ranging from 3 to 1,000 hours. In the future, we plan to extend FHVAE applications, such as ASR domain adaptation [7] and audio conversion [5] to larger scales.

<sup>2</sup>In these images, harmonic spacing is the clearest cue for fundamental frequency differences. Far-field recordings tend to have lower signal-to-noise ratios, which results in blurrier images.

<sup>3</sup>Phonetic content can usually be determined by the spectral envelope, and relative position of formants.

## 7. References

- [1] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2011.
- [2] S. Tan and K. C. Sim, "Learning utterance-level normalisation using variational autoencoders for robust automatic speech recognition," in *Spoken Language Technology Workshop (SLT), 2016 IEEE*. IEEE, 2016, pp. 43–49.
- [3] J. F. Drexler, "Deep unsupervised learning from speech," 2016.
- [4] W.-N. Hsu, Y. Zhang, and J. Glass, "Learning latent representations for speech generation and transformation," in *Interspeech*, 2017, pp. 1273–1277.
- [5] —, "Unsupervised learning of disentangled and interpretable representations from sequential data," in *Advances in Neural Information Processing Systems*, 2017.
- [6] W.-N. Hsu and J. Glass, "Extracting domain invariant features by unsupervised learning for robust automatic speech recognition," in *Acoustics, Speech and Signal Processing (ICASSP), 2018 IEEE International Conference on*. IEEE, 2018.
- [7] W.-N. Hsu, Y. Zhang, and J. Glass, "Unsupervised domain adaptation for robust speech recognition via variational autoencoder-based data augmentation," in *Automatic Speech Recognition and Understanding (ASRU), 2017 IEEE Workshop on*. IEEE, 2017.
- [8] L. Van Der Maaten, "Accelerating t-sne using tree-based algorithms," *Journal of machine learning research*, vol. 15, no. 1, pp. 3221–3245, 2014.
- [9] J. S. Garofalo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, and D. S. Pallett, "DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. NIST speech disc 1-1.1," *NASA STI/Recon technical report n*, vol. 93, 1993.
- [10] D. Pearce, "Aurora working group: Dsr front end lvcsr evaluation au/384/02," Ph.D. dissertation, Mississippi State University, 2002.
- [11] J. Carletta, "Unleashing the killer corpus: experiences in creating the multi-everything ami meeting corpus," *Language Resources and Evaluation*, vol. 41, no. 2, pp. 181–190, 2007.
- [12] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 5206–5210.
- [13] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlicek, Y. Qian, P. Schwarz *et al.*, "The kaldi speech recognition toolkit," in *IEEE 2011 workshop on automatic speech recognition and understanding*, no. EPFL-CONF-192584. IEEE Signal Processing Society, 2011.
- [14] J. Garofalo, D. Graff, D. Paul, and D. Pallett, "Csr-i (wsj0) complete," *Linguistic Data Consortium, Philadelphia*, 2007.
- [15] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [16] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [17] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, "Tensorflow: A system for large-scale machine learning," in *OSDI*, vol. 16, 2016, pp. 265–283.
- [18] A. K. Halberstadt, "Heterogeneous acoustic measurements and multiple classifiers for speech recognition," Ph.D. dissertation, Massachusetts Institute of Technology, 1999.